



University of Padua  
Department of Information Engineering



# The University of Padua at CLEF 2004: Experiments on Statistical Approaches to Compensate for Limited Linguistic Resources



Giorgio Maria Di Nunzio - Nicola Ferro - Nicola Orio

{dinunzio, nf76, orio}@dei.unipd.it

Workshop of the Cross Language Evaluation Forum 2004 (CLEF 2004)  
Bath, UK, 15-17 September 2004

# History of the IMS at CLEF

---

## – 2002 Monolingual

- Language independent stemmer
- IRON

## – 2003 Mono/Bilingual

- Probabilistic models for automatic stemmer generation
- Web IRON

## – 2004 Mono/Bilingual

- Limited Language Resources
  - For stemming and query translation
- IRON enhanced

# Main Objectives

---

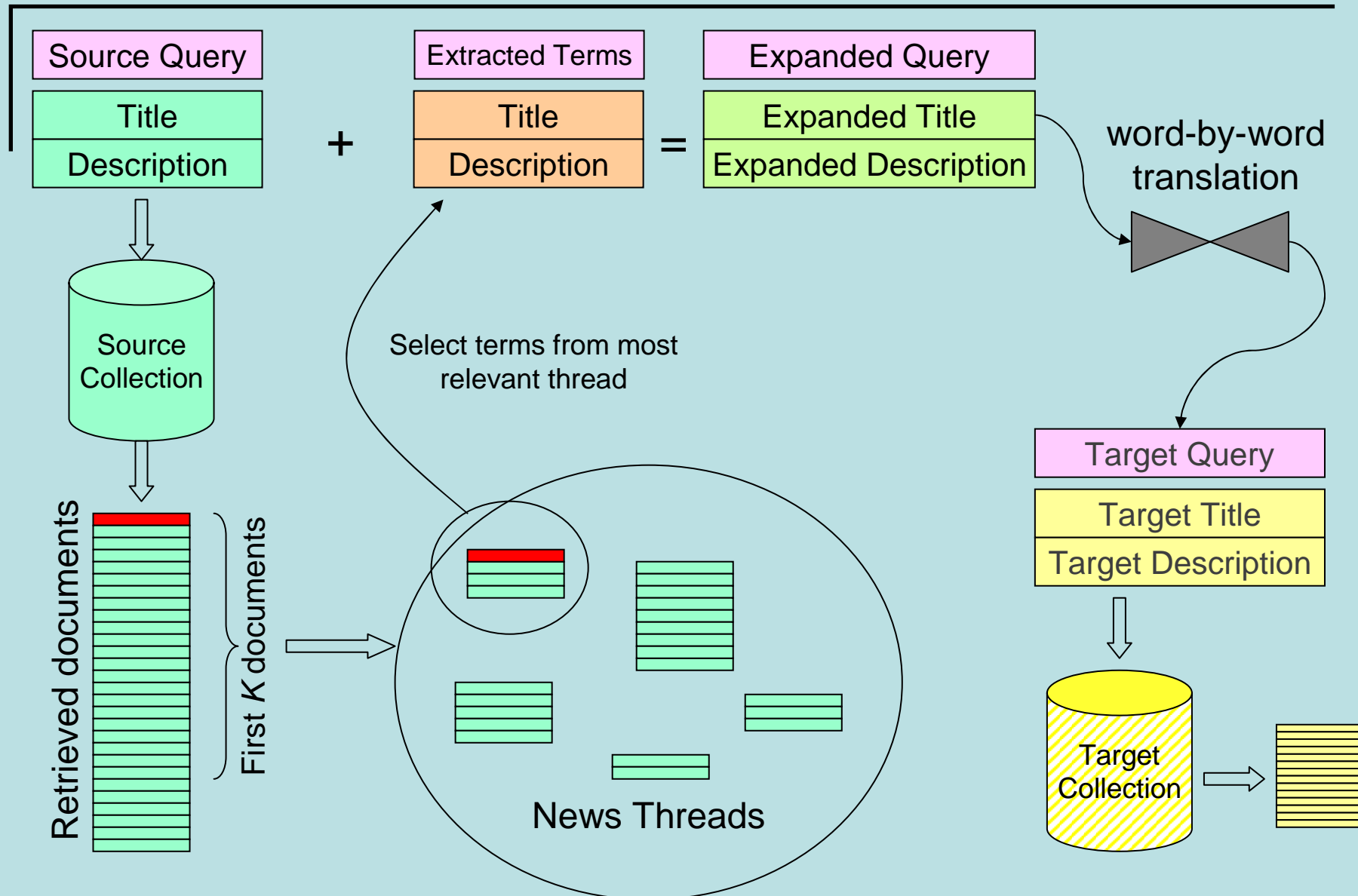
- Minimize human efforts when applying IR techniques to new languages
- Partially overcome problems of limited language resources
  - Lack of advanced tools for query/document translation
  - Possible lack of knowledge on morphological structure for stemming
- Improve our evaluation prototype system

# Bilingual

---

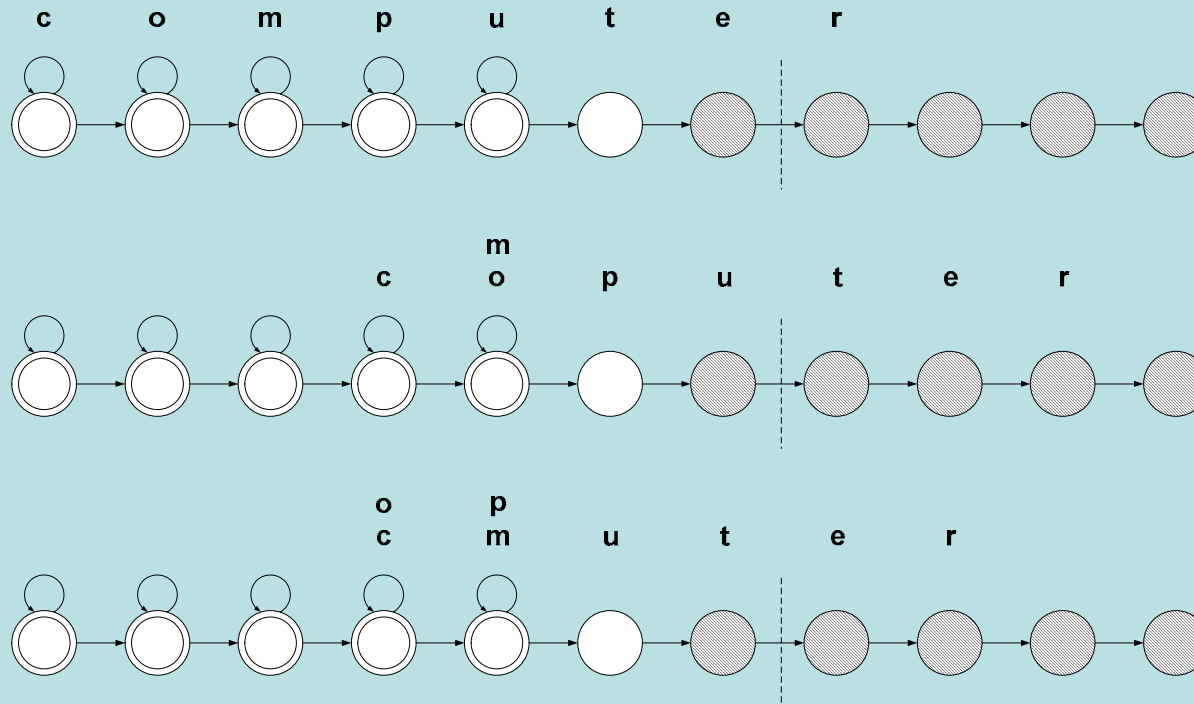
- Almost Comparable Corpora
  - *Automatic news thread identification* by means of hierarchical clustering
- Query expansion
  - Extract significant terms from the most relevant news thread
- Query translation
  - Translate expanded query using on-line word-by-word translation services (Google)
  - No control on the size of the vocabulary
  - No synonyms available

# Bilingual

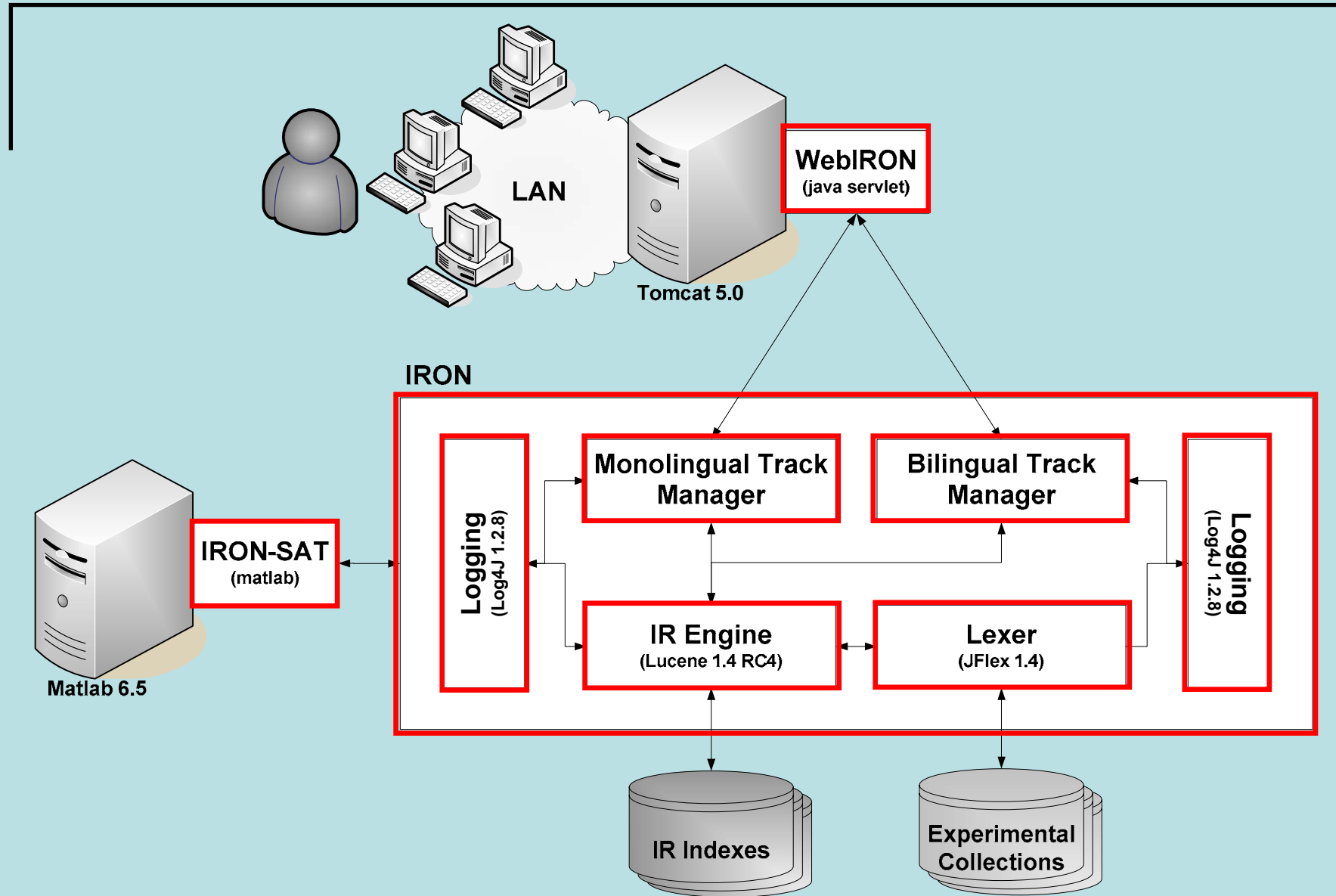


# Monolingual

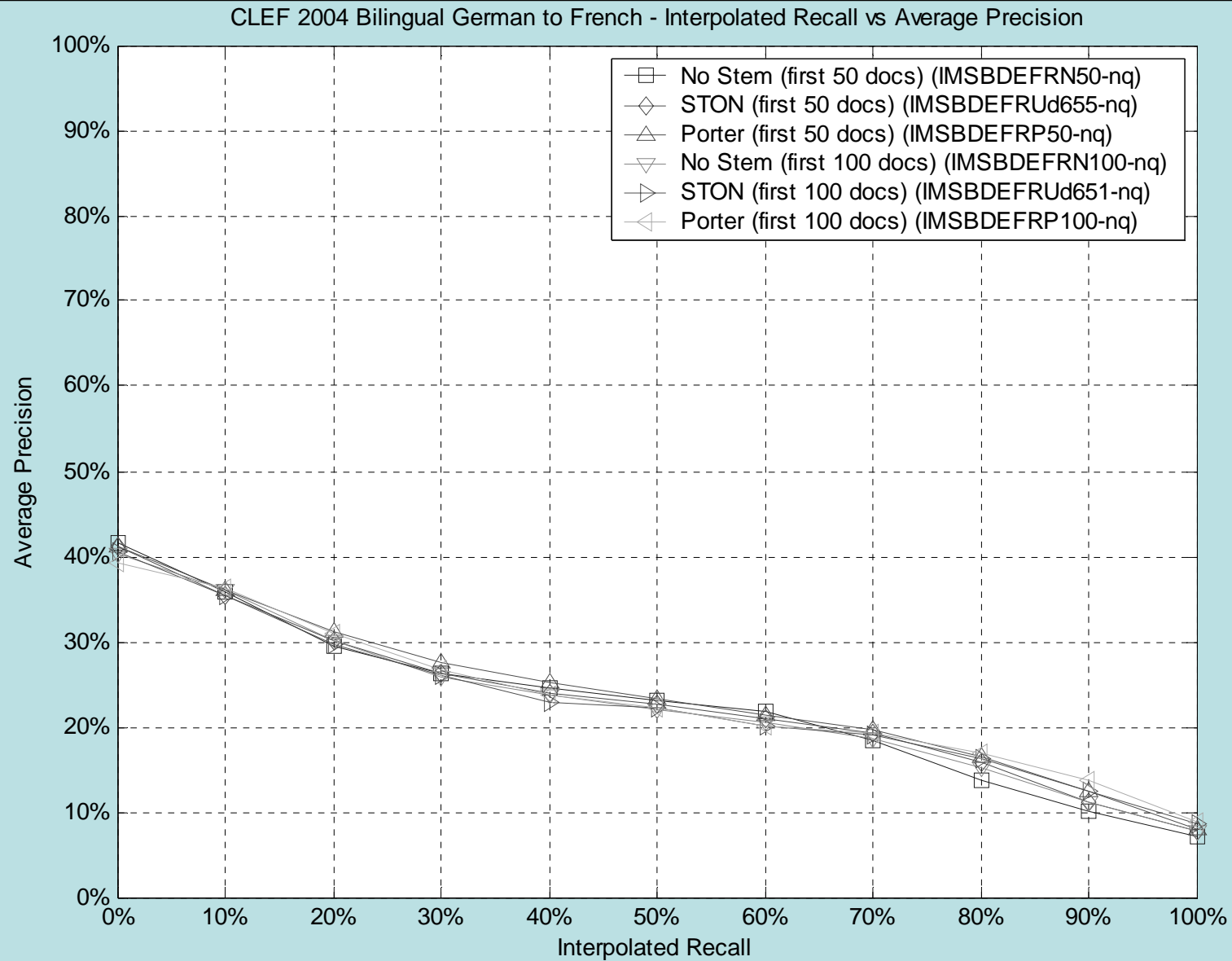
- STON : Hidden Markov Models stemmer
  - The sequence of letters of a word can be considered as a sequence of symbols emitted by a HMM
  - Most probable path for the observed word
  - Transition from stem-set to suffix-set (split-point)
  - Only needs a set of words of the language to be trained off-line



# IRON System

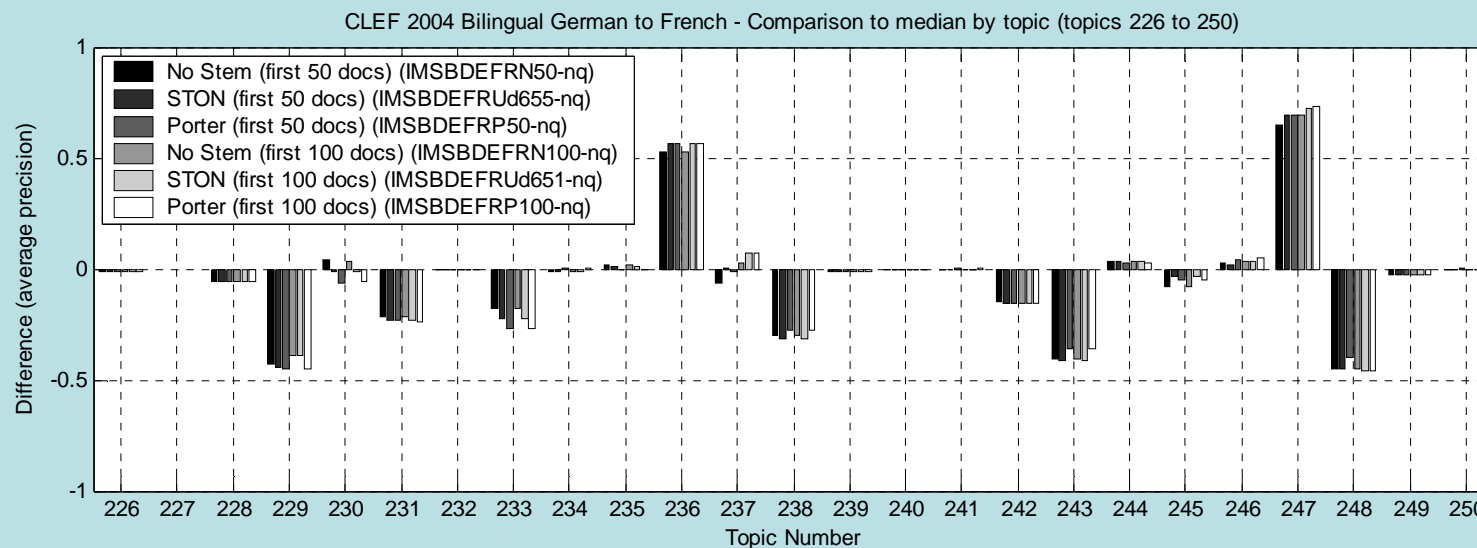
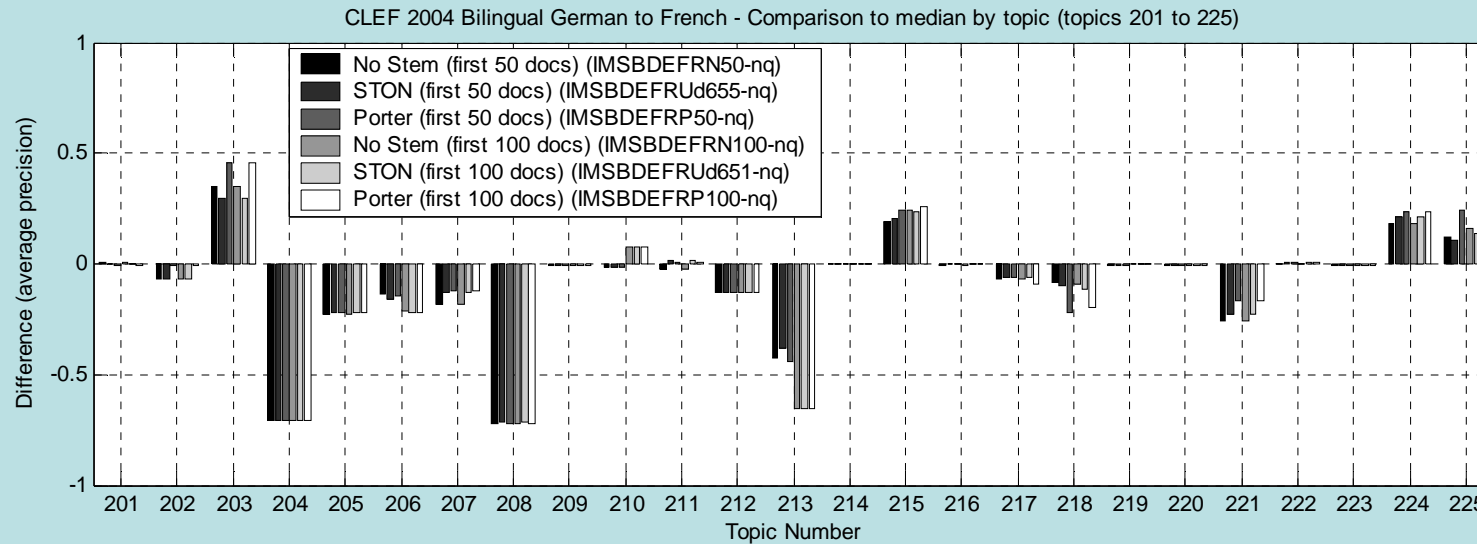


# Bilingual Experiments

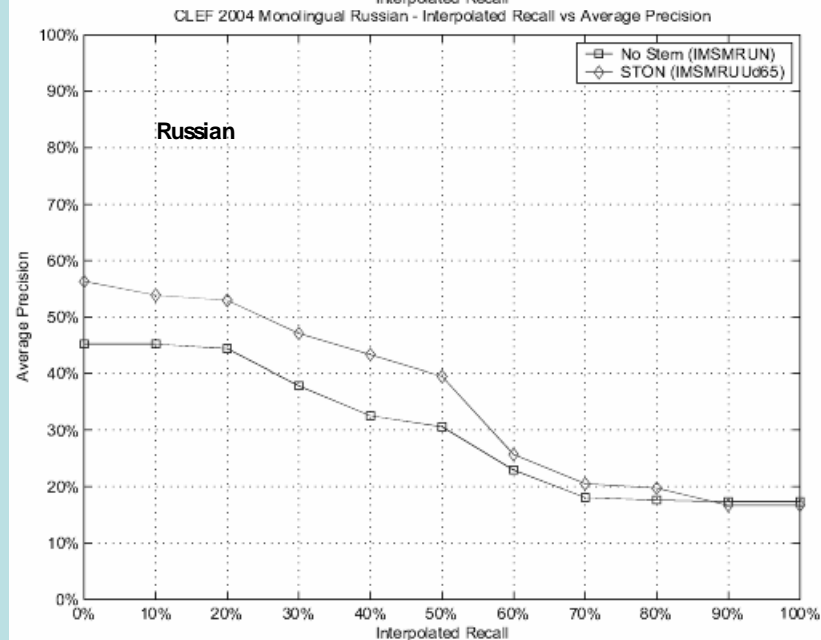
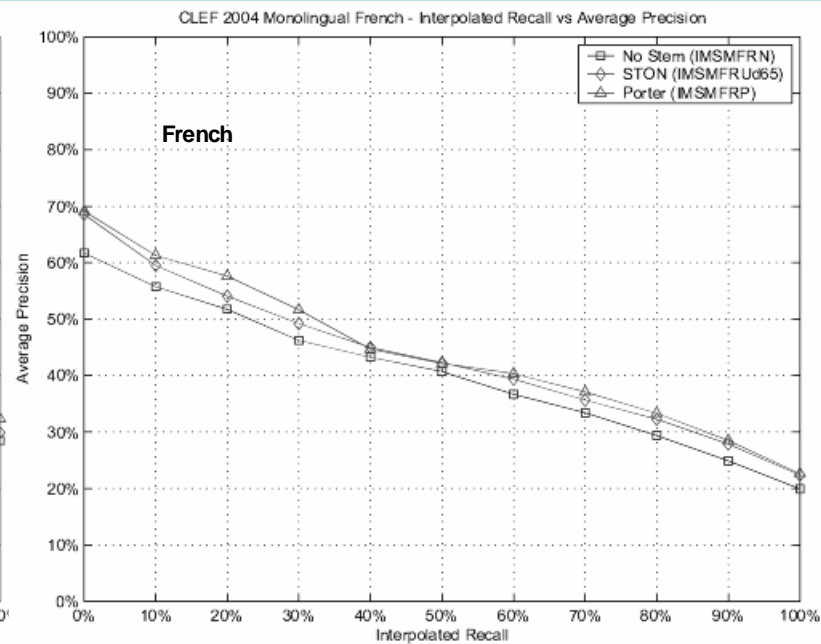
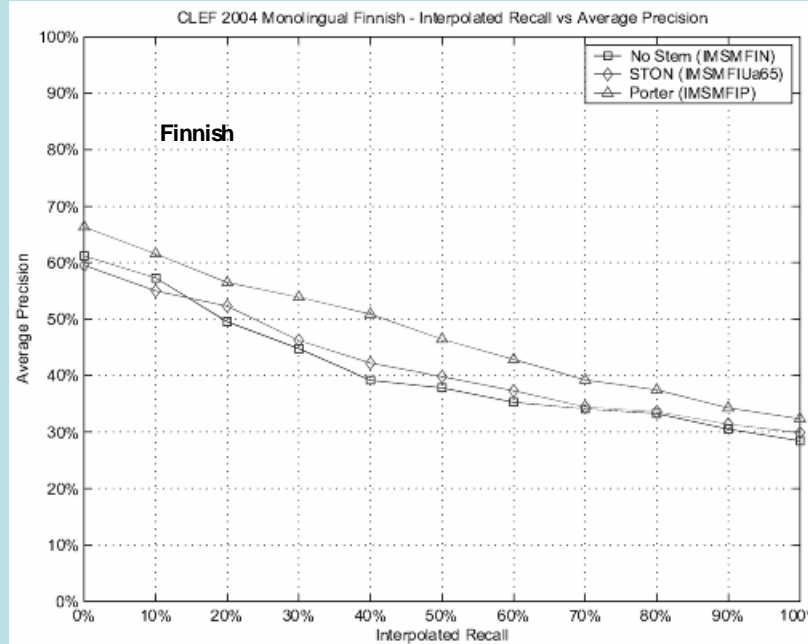




# Bilingual Experiments



# Monolingual Experiments



Comparison STON - No Stem (relative %)

	Finnish	French	Russian
<b>Rel. Retr.</b>	1.93 %	0.33 %	0.39 %
<b>Avg. Prec.</b>	65.83 %	29.15 %	6.79 %
<b>Exact R-Prec.</b>	50.28 %	20.85 %	57.81 %

# Conclusions and Future Work

---

- Minimizing human labor and language resources
    - Automatic stemmer generation
    - Automatic query expansion by means of hierarchical clustering
    - Free on-line word-to-word translation
  - Statistical analysis of results
- 
- Thread identification in both source and target collection.
    - Coupling between threads to refine results